

Compressing graphs with semantic structure

Elaine Angelino

Abstract. Graph compression schemes are based on finding an ordering or clustering of nodes that places similar nodes close to one another. Existing algorithms for Web graphs and social networks tend to only consider nodes and edges, where each node comes with some identifier (a name or number), and edges are either directed or undirected but otherwise unlabeled. However, real world graphs tend to come with additional semantic information that is typically ignored by the compression algorithms. For example, nodes corresponding to individuals in a social network are usually associated with structured user profiles. Here, we begin to study whether this kind of semantic structure can be leveraged to produce compression-friendly orderings. We emphasize that this additional semantic structure comes “for free” in the sense that it is usually already collected and stored by the curators of the graph datasets. Our preliminary results show that semantic ordering heuristics are useful and may be competitive with other orderings for social networks.

1. INTRODUCTION As large graph datasets are collected and curated in more domains, there is an increasing need for efficient graph storage. Web graphs, social networks, and protein interaction networks are just some examples of large graphs corresponding to real world data, with hundreds of thousands to millions or billions of nodes, and millions to billions of edges. These nodes and edges are often accompanied by rich semantic information that ranges from structured data (e.g. additional node or edge attributes) to unstructured text and other data (e.g. Web page content).

While most graph compression schemes consider the “core” graph data – i.e. the nodes and edges of a graph – they tend to ignore the other semantic data. This is unfortunate for several reasons. (1) Graph data curators collect, store and analyze the graph data in the context of extra semantic data. For example, social networking and commercial sites probably want to collect as much user information as possible, for business analytics and marketing applications. Meaningful compression in real systems requires performance over core graph data plus extra semantic data. (2) Intuition tells us that semantic features may correlate with graph properties. For example, nodes with similar semantic information may be more likely to have similar neighborhood relationships. Especially since many graph datasets already include extra semantic information, it would be attractive if this existing structure could be leveraged for graph compression and other related applications, e.g. clique-finding, community detection, etc.

Graph compression schemes are based on finding an ordering or clustering of nodes that places nodes with similar neighborhood relationships close to one another. For example, a graph that decomposes into cliques plus just a few edges between cliques can be encoded as lists of nodes, one for each clique, plus a list of the extra edges. This can be a much more efficient encoding than enumerating all edges, but clique-finding is well-known to be a hard problem. Likewise, computing an optimal ordering of nodes has been shown to be NP hard [Chierichetti et al. 2009].

Several efficient ordering heuristics provide impressive compression performance. For example, the BV compression scheme is based on the fact that simply lexicographically sorting nodes in a Web graph by URL provides a natural, compression-friendly ordering [Baldi and Vigna 2004]. This works well because of lexicographic locality: (1) pages that are close in this ordering tend to have similar neighborhood relationships, and (2) many links occur within the same domain, and will appear near one another under this ordering. A simple social network, where nodes correspond to individuals and edges correspond to social links (e.g. ‘is friends with’, ‘works with’), does not by itself suggest a natural ordering analogous to the lexicographic ordering for Web graphs. This motivates the shingle ordering heuristic, which is efficient and yields good compression performance [Chierichetti et al. 2009].

However, as we pointed out, most graph datasets are actually more than just nodes and edges. In a social network, each node has a unique identifier; these identifiers do not necessarily provide a natural compression-friendly ordering. At the same time, each node is typically associated with semantically rich additional data, e.g. a user profile containing data for fields such as ‘last name’, ‘first name’, ‘geographic location’, ‘profession’,

'employee', 'college', 'interests', etc. This structured information can be thought of as a vector of features associated with each node. Intuition tells us that nodes may tend to have neighbors with similar feature vectors, at least over some subset of fields. This idea plays an important role in many applications, such as various clustering and community detection algorithms for social networks, user preference prediction (e.g. the motivation for the Netflix challenge), and coarse heuristics for sharding large graph datasets across multiple servers (e.g. partitioning user accounts by geographic region).

We are generally interested in exploiting structured semantic information to compress graphs. In this paper, we evaluate the usefulness of rudimentary semantic heuristics for obtaining compression-friendly orderings on nodes. Our basic idea is to obtain node orderings by lexicographically or numerically sorting associated structured semantic information, e.g. sort nodes by a user-provided tag, or by tag followed by timestamp. For the dataset tested, we observe compression performance under semantic ordering heuristics to be competitive with shingle ordering heuristics with a single shingle, and potentially comparable to double shingle orderings. The preliminary results encourage further evaluation with additional graph datasets and more sophisticated semantic ordering heuristics.

2. RELATED WORK Graph compression techniques have thus been developed in several domains, including for Web graphs [Boldi and Vigna 2004] and social networks [Chierichetti et al. 2009]. A general approach in graph compression is to determine an ordering or clustering of nodes that places similar nodes close to one another; a node N_i can then be efficiently represented in terms of differences with respect to node N_{i-1} . Determining the best ordering is NP-hard, but "natural" and heuristic orderings have demonstrated good results [Chierichetti et al. 2009].

Compression techniques for the Web graph provide a foundation for compressing social networks. Web graphs are highly compressible due to their property of "lexicographic locality"; we described this idea in § 1, and it forms the basis of the Boldi-Vigna (BV) compression scheme [Boldi and Vigna 2004]. Social networks lack a natural ordering, but a shingling heuristic places nodes with similar neighbors close to one another. This provides a basis for a BV-like scheme called backlinks (BL), that further takes advantage of the reciprocal nature of social networks [Chierichetti et al. 2009].

If the goal is to study graph compression for application to real systems, then it is important to design compressed data structures that meet the access and computational needs of real systems. This is the motivation for the recently developed layered label propagation (LLP) method, which produces a compressed graph data structure that corresponds to a hierarchical clustering of nodes [Boldi et al. 2010]. While many graph clustering algorithms tend to produce one particularly large node, the LLP scheme is designed to avoid this problem. The hierarchical nature of the compressed data structure is particularly attractive for systems

+We note that real storage systems for managing graph data are still in their infancy compared to relational database management systems, and many do not exploit the graphical structure of their data. In fact, the systems tend to map graph data onto a relational framework. For example, the RDF-3X engine treats a graph as a single relation of (node, edge, node) triples, and stores these in a compressed clustered Btree [Neumann and Weikum 2008]. The triples are sorted lexicographically and compressed via dictionary encoding. RDF-3X depends on an exhaustive set of indexes on top of this structure to achieve fast query performance.

There is a rich literature of graph partitioning and clustering algorithms, which can be applied to graph compression because they tend to group similar nodes, e.g. [Fe der and Motwani 1995]. A thorough review is beyond the scope of this paper, but we note relevant work in community detection [Flake et al. 2000; Iino et al. 2005] and modularity [Newman 2003; 2006].

3. COMPRESSION SCHEME Our compression scheme is very similar to the backlinks (BL) compression scheme for social networks

[Chierichetti et al. 2009]. The are two differences that are mainly due to our choice of dataset: (1) instead of a social network, we consider a bipartite graph (e.g. two node types correspond to documents and words, and edges connect documents to the words they contain), and (2) our edges are undirected, and so we have no notion of link reciprocity. By using ideas from the BL compression scheme – in particular, link reciprocity – we could easily extend our scheme to encode directed bipartite graphs.

Consider a bipartite graph $G(N_1, N_2, E)$ where N_1 and N_2 are the nodes of the two distinct types, $|N_1|$ and $|N_2|$ are the number of nodes in each set, and E is the set of edges. Let D be the largest difference, and let $k = \lfloor \log_2 D \rfloor$.

Let v_1, \dots, v_n be the nodes in N_1 and v_{n+1}, \dots, v_{n-1} be the nodes in N_2 . Let E be the set of edges between N_1 and N_2 .

Let u_1, \dots, u_n be the nodes in N_1 and u_{n+1}, \dots, u_{n-1} be the nodes in N_2 . Let E be the set of edges between N_1 and N_2 .

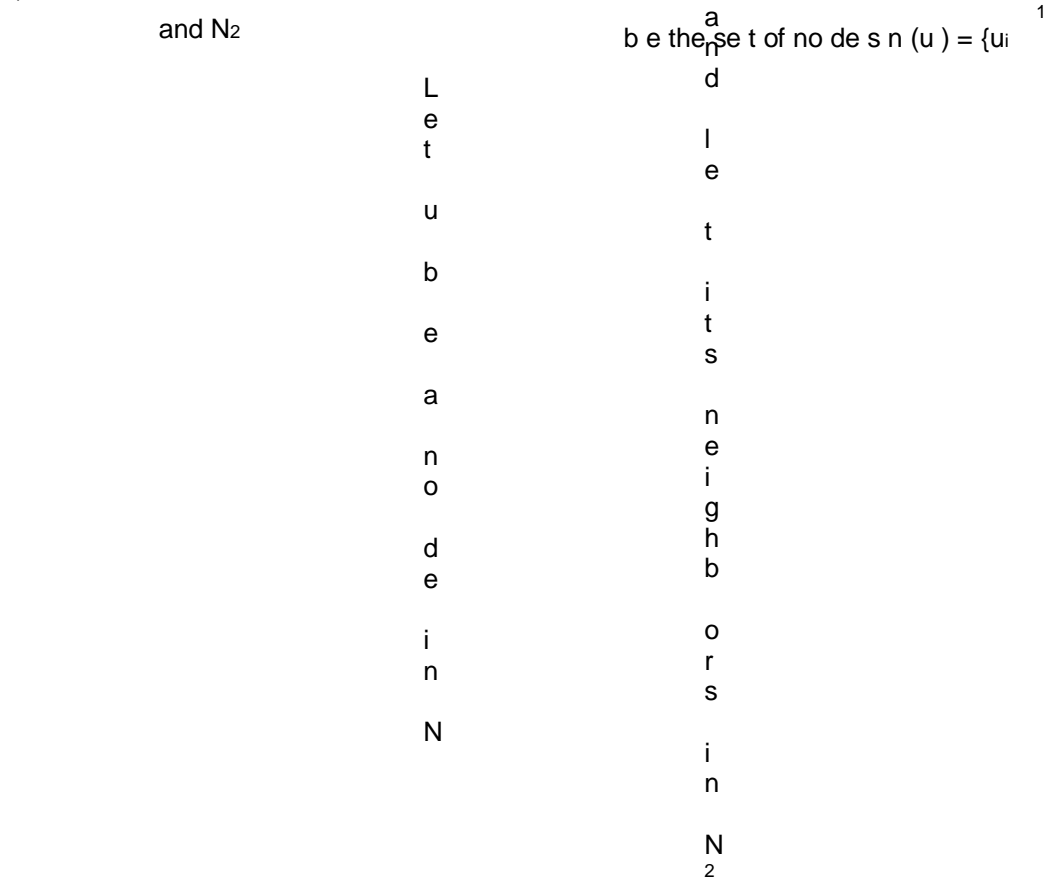
Let u_1, \dots, u_n be the nodes in N_1 and u_{n+1}, \dots, u_{n-1} be the nodes in N_2 . Let E be the set of edges between N_1 and N_2 .

we have dictionary mapped the nodes in N_1 to integers in $\{1, \dots, |N_2|\}$. Let v be a node in N_1 . We encode v with the following pieces of information:

- (1) Base information. The degree of v , $\deg(v)$. This requires at most $m = \lceil \log_2(|N_2| + 1) \rceil$ bits.
- (2) Copying. The preceding node u that v uses as a prototype. We encode u as the difference between u and v . If $u = v$, we copy nothing. Otherwise, we write down $\deg(u)$ bits, each representing whether a neighbor of u is also a neighbor of v .
- (3) Residual edges. Now we must encode the remaining neighbors of v not described by the copying step. Let the set of neighbors be $v_1 < \dots < v_n$. We write down v and encode the gaps $|v(D + 1)|$. Writing down k requires $\lceil \log_2(k + 1) \rceil$ bits. Decoding the residual edges then requires $m + (n - 1) \cdot k$ bits.

4. NODE ORDERINGS As in BV or BL compression, the performance of our compression scheme depends on the ordering of the nodes. While finding an optimal ordering for social networks is NP-hard, a shingling heuristic gives good results and is easy to compute [Chierichetti et al. 2009]. We describe two compression-friendly orderings: (1) a straightforward extension of the shingling heuristic for bipartite graphs, and (2) a simple ordering heuristic that leverages structural semantic information associated with the nodes in the graph.

We use the same notation as in § 3. Consider a bipartite graph $G = (N_1, N_2, E)$,



$\{s(u_i) \mid u_i \in N_1\}$
The nodes in N_1

$1, N$

1
 2

1

The nodes in N_1

1

are the nodes of the two distinct types, $|N|$, and the edges in E are undirected and represent links between nodes in N . We want to find a compression-friendly ordering of the nodes in N . 4.1 Shingle ordering heuristic

Let s be a random permutation of the nodes in N . The shingle of u is the minimum element in $n(u)$ under s , $M(u) = s$. In practice, a min-wise independent family or even pairwise independent hash functions can be used for s [Broder et al. 1998; Chierichetti et al. 2009].

can be ordered by the ir shingles. If two nodes in N share many of the same neighbors in N , they will with high probability be near each other in the shingle ordering. Notice that the shingle ordering gives a partial ordering of the nodes in N . A shingle ordering can be refined by combining shingles obtained under independent hash functions: the ordering under a second shingling can be used to break ties in the first, etc.

4.2 Semantic ordering heuristic Now, we consider additional structured semantic information about our graph. While it would be very

interesting to consider semantic information on edges, we restrict ourselves here to semantic information on nodes. Let $G(N, E; S)$ be our bipartite graph as before, plus semantic information S associated with each node in N . We can think of an element in S as a d -tuple whose k th element corresponds to the k th feature in a space of d features.

can be ordered according to a sort with respect to some projection of the associated features in S . This may involve lexicographic and/or numerical sorting, depending on the nature of the semantic features. For example, suppose S has three fields corresponding to 'name' (string), 'time' (integer representing seconds since the epoch), and 'tag' (string). Lexicographically sorting the 'name' field provides a 1-dimensional ordering of the nodes in N . Lexicographically sorting the 'tag' field followed by numerically sorting the 'time' field provides a 2-dimensional ordering of the nodes in N . These orderings based on multi-dimensional projections of our feature space are similar in spirit to the method of combining shingles described above. We note that the separate dimensions in our feature vector are in general not independent.

3
 N

Equivalently, we could encode $G(N$

, E) plus the semantic information S as a $(d + 2)$ -partite graph. This may be a useful representation worth further consideration.

5. EXPERIMENTS Our primary goal is to compare semantic ordering heuristics to shingle ordering heuristics. We consider semantic orderings based on individual as well as multiple semantic fields, and shingle orderings based on individual and multiple shingles. In this section, we describe the dataset we use for our experiments, baseline orderings for comparison, and some implementation details.

5.1 Data We focus on one particular dataset, a set of tweets from June 2009 obtained by crawling the popular microblogging website, twitter.com [Yang and Leskovec 2010].² Each tweet is a short message consisting of up to 140 characters, and is associated with two additional pieces of data: a username and a time stamp (with second-level granularity). Each tweet itself may contain two particularly interesting kinds of user-provided information: hashtags and the use of the '@' symbol to denote a reply. A hashtag is a token (e.g. a word) that starts with the '#' symbol; users can group and search for tweets by hashtag, and so are motivated to provide semantically meaningful hashtags. The '@' symbol followed by a username indicates that the tweet is a reply to the user associated with username.

We think of this data as a bipartite graph $G(T, V, E; S)$ where T is the set of nodes corresponding to tweets and V is the set of nodes corresponding to tokens in the tweet vocabulary. An edge between $t \in T$ and $v \in V$ indicates that v is a token in tweet t ; the edges are undirected. This bipartite graph represents a "bag of words" interpretation of the corpus of tweets. We define the additional semantic features S on the tweets to include the associated username and timestamp, as well as any tokens that can be interpreted as hashtags or replies.

We summarize the basic properties of our dataset in Table I (a). For simplicity, we eliminate all tweets with Unicode characters. This is a crude way to get rid of tweets, such as those written in Asian characters, that require different lexicons for tokenization. Namely, we only consider tweets that can be lexed by splitting the tweet string on white space; our method happens to eliminate more tweets than is strictly necessary. We also only consider tweets containing at least one token that looks like a hashtag (we do not consider the '#' symbol alone to be a hashtag). This yields about 1.7 million tweets (about 10% of the original dataset) and 2.0 million tokens connected by 24 million edges.

We summarize the semantic features of our dataset in Table I(b), i.e. the number of unique users, timestamps, extracted replies and hashtags. In the case of multiple replies, we extract only the first. In the case of multiple hashtags, we extract only the first two, and consider each as a separate feature.

(a) Basic graph properties.		(b) Semantic features of tweets.	
Notation	Feature Number	Feature Number	
T tweets	1,731,402	V tokens	2,073,037
E links	24,042,218	edges	24,042,218
		primary	5,569
		secondary	13,854
		replies	1,731,072
		hashtags	53,384
		users	423,616
		timestamps	83

Table I. Summary of (a) basic graph properties and (b) semantic features of our Twitter dataset.

5.2 Baseline orderings For comparison purposes, we consider two baseline orderings of the nodes in T :

- (1) Random order. This is a random permutation of the nodes.
- (2) Natural order. This is the order of the nodes in the original dataset, presumably corresponding to the order in which the twitter.com pages were crawled. We note that orderings based on individual pieces of semantic information are often considered as a baseline. For example, zip codes of individuals in a social network can be used to obtain a geographic ordering [Chiericetti et al. 2009].

²We note that this data is apparently no longer available through the SNAP (Stanford Network Analysis Platform) website, "a requester from Twitter."

5.3 Implementation details (1) Compression scheme. Given a node ordering, to encode a particular node u , we only consider the preceding node in the ordering for the prototype of u , i.e. we only use a window size of 1. (2) Hash functions. We use NumPy's `numpy.random.permutation()` function to generate permutations of $\{1, \dots, |V|\}$ for our shingle ordering heuristic.

6. RESULTS In this section, we summarize our compression performance results, comparing shingle order heuristics to semantic order heuristics, for single orders, double orders, and additional orders (Table I and Figure 1). We also include results for our baseline orderings. All compression results are given in terms of bits per link.

(a) Single orders. Ordering Bits / link	(b) Semantic double orders. hashtag time reply user	(c) 3- and 4-orders. Ordering Bits / link																																	
shingle 17 .455* hashtag ag 17.629 user 18.700 reply 19.142 time 20.334 natural 20 .331 random 21.042	<table border="1"> <tr> <td>hashtag - 16</td> <td>.513</td> <td>16.637</td> <td>16.874</td> <td>time</td> </tr> <tr> <td>20.270 - 20</td> <td>306</td> <td>20.327</td> <td>reply 17</td> <td>.011 18.220 - 17</td> </tr> <tr> <td>.707</td> <td>user 18</td> <td>.141</td> <td>18.026</td> <td>18.249 -</td> </tr> </table>	hashtag - 16	.513	16.637	16.874	time	20.270 - 20	306	20.327	reply 17	.011 18.220 - 17	.707	user 18	.141	18.026	18.249 -	<table border="1"> <tr> <td>3-shingle 15</td> <td>.135</td> <td>hashtag, time,</td> </tr> <tr> <td>reply 16</td> <td>.501</td> <td>hashtag, time, user 16</td> </tr> <tr> <td>.499</td> <td></td> <td></td> </tr> <tr> <td>4-shingle 14</td> <td>.112</td> <td>hashtag, time,</td> </tr> <tr> <td>reply, user 16</td> <td>.498</td> <td>hashtag, time,</td> </tr> <tr> <td>user, reply 16</td> <td>.499</td> <td></td> </tr> </table>	3-shingle 15	.135	hashtag, time,	reply 16	.501	hashtag, time, user 16	.499			4-shingle 14	.112	hashtag, time,	reply, user 16	.498	hashtag, time,	user, reply 16	.499	
hashtag - 16	.513	16.637	16.874	time																															
20.270 - 20	306	20.327	reply 17	.011 18.220 - 17																															
.707	user 18	.141	18.026	18.249 -																															
3-shingle 15	.135	hashtag, time,																																	
reply 16	.501	hashtag, time, user 16																																	
.499																																			
4-shingle 14	.112	hashtag, time,																																	
reply, user 16	.498	hashtag, time,																																	
user, reply 16	.499																																		

Table I. Summary of compression performance for different orderings and combinations of orderings. All results are reported in bits per link. Note that in these tables, "hashtag" refers to the primary hashtag, i.e. the first hashtag appearing in a tweet. (a) Results for single orderings. * This result for the single shingle ordering represents the best from 10 independent trials (mean = 17.690 bits/link, standard deviation = 0.186 bits/link, max = 18.172 bits/link). (b) Results for double semantic orderings. The rows and columns correspond to the first and second semantic features used in the ordering. Note that the double shingle ordering beats all of these by at least 1 bit/link, e.g. achieving 15.135 bits/link with two particular randomly chosen shingles. Also note that the ordering achieved by combining primary and secondary hashtag data gives 17.075 bits/link. (c) Orderings that combine three or four shingles or semantic features.

6.1 Baseline orderings For comparison purposes, we report results for our two baselines in Table I(a). Compression with the natural and random orders require about 20.3 and 21.0 bits/link, respectively.

6.2 Shingle orderings

We generate 10 different shinglings. First, we consider each shingle separately. Over the 10 trials, compression performance ranged from about 17.5 - 18.2 bits/link, with a mean and standard deviation of 17.7 and 0.2 bits/link, respectively. Note that the best performance over the 10 trials is reported in Table I(a).

We also consider combinations of increasing numbers of shingles, from 1 to 10, for a particular randomly chosen order of the shingles. We note that we can obtain significant performance gains by combining multiple shingles, but that the gains decay quickly as additional shingles are included (Figure 1). In our particular example, the double shingle ordering gives an improvement of about 2.5 bits/link. Compression with the 4-shingle ordering requires only 13.847 bits/link, the 5-shingle ordering decreases this by 0.096 bits/link, and the 10-shingle ordering decreases by an additional 0.094 bits/link. We also report these results for 3- and 4-shingle orderings in Table I(c).

6.3 Semantic orderings First, we consider orders obtained by sorting along the following single semantic dimensions: primary hashtag, user, time stamp, reply (Table I(a)). The ordering by time stamp gives compression comparable to the natural order; it turns out the tweets we're already nearly ordered by time stamp. The other three single semantic orderings are all better than the natural ordering. In particular, the hashtag ordering is the best, and is in fact comparable to the single shingle orderings.

Next, we consider all 12 orders obtained by sorting according to ordered pairs chosen from the four se

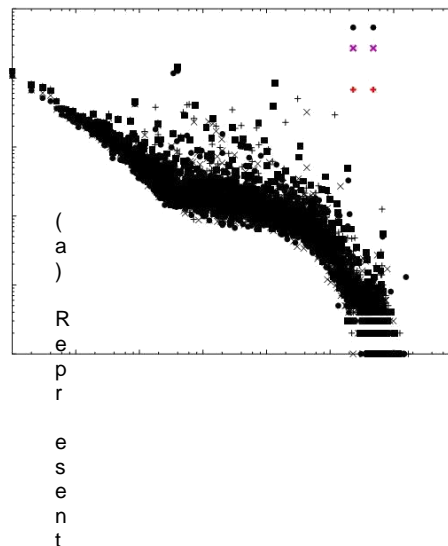
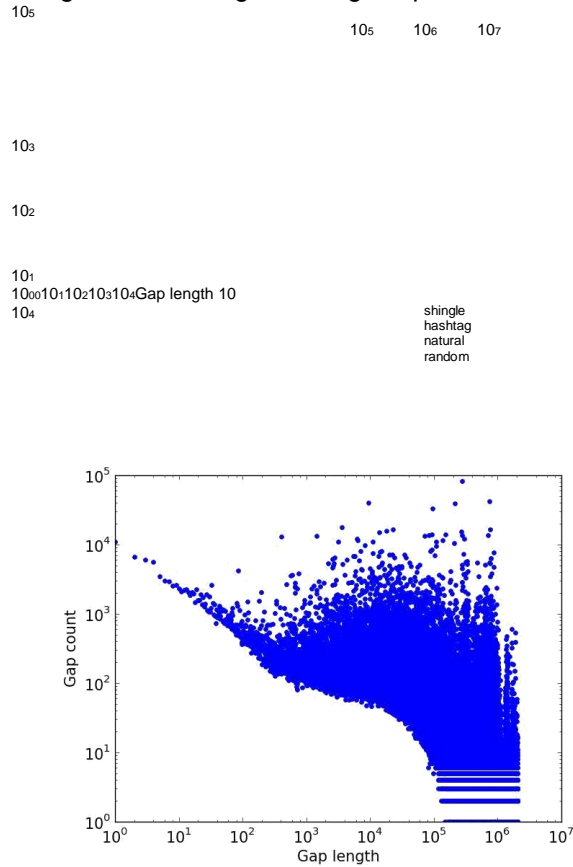
mantic dimensions above (Table I I(b)). We also tested compression obtained by None of the double semantic orderings are as good as the double shingle ordering (about 15 bits/link), with the best – ordering

13.5 14.0 14.5 15.0 15.5 16.0 16.5 17.0 17.5 18.0 19.0
 10 Number of shingles

Fig. 1. Multi-shingle orderings. Bits per link as a function of number of shingles.

by hashtag followed by timestamp – falling behind by about 1.5 bits/link. Except for double semantic orderings that start sort first by timestamp, addition of a second semantic ordering to a first improves performance by about 0.5- 2.0 bits /link, which is less than the 2.5 bits/link gain we obtained by using the double shingle compared to the single shingle ordering. Combinations of three to four semantic orderings only give small incremental improvement, some examples are shown in Table I(c).

6.4 Gap lengths and residual edges Chierichetti et al argue that the shingle ordering enables good compression by creating more small gaps, compared to the other orderings. They experimentally show this by plotting the number of gaps versus gap length for different orderings, which exhibit very different distributions [Chierichetti et al. 2009]. We produce similar plots in Figure 2, but observe only very subtle differences between the distributions for the different orderings. Figure 2(a) plots a representative distribution, the number of gaps (y-axis) versus gap length (x-axis) for the hashtag ordering. Following Chierichetti et al, Figure 2(b) plots a sub-sample of the distributions, for shingle (blue circles), hashtag (magenta x's), natural (cyan squares), and random (red crosses) orderings. In particular, we were surprised that there was not an immediate difference between the distributions for the shingle (or hashtag) versus random (or natural) orderings. We expected shingle and hashtag orderings to produce more small gaps than the random and natural orderings.



ative distribution).

Bits per link

(
b
)
S
u
b
-
s
a
m
p
l
e
d
g
a
p
d
i
s
t
r
i
b
u
t
i
o
n
s
.

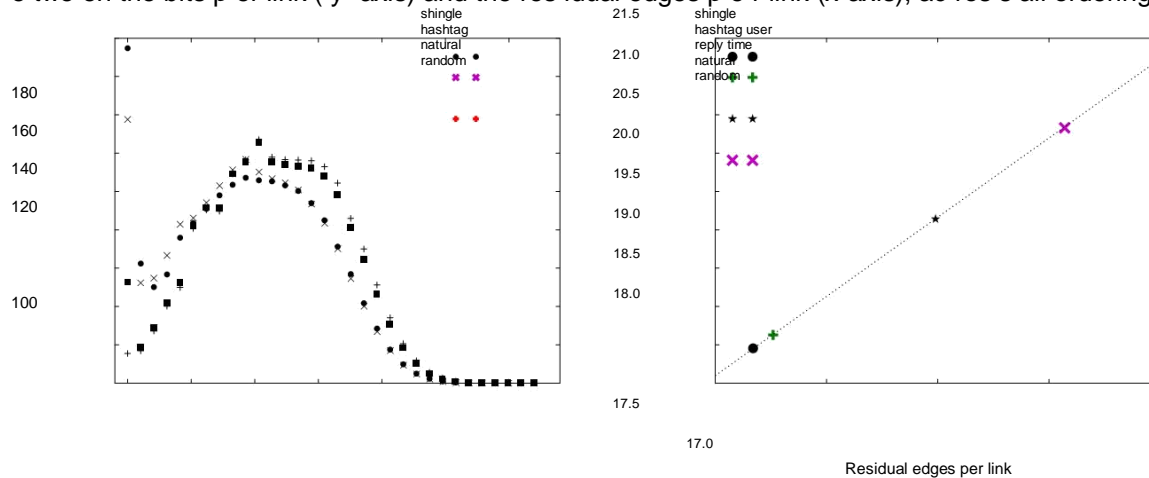
Fig. 2 . Number of gaps (y -axis) of a certain length (x -axis) , for (a) a representative ordering (hashtag) , and (b) a sub-sampling of points for shine (blue circles) , hashtag (magenta x's) , natural (cyan squares) , and random (red crosses) orderings.

6

Gap count

Since the distribution of gap lengths did not explain the differences between the orderings, we decided to look at the number of residual edges in the encoding, for the different orderings. Figure 3(a) plots the number of tweets (in thousands, y-axis) with a certain number of residual edges (x-axis), for shingle (blue circles), hashtag (magenta x's), natural (cyan squares), and random (red crosses) orderings. Here, we can see that the shingle and hashtag orderings allow more tweets to be encoded with a smaller number of residual edges, i.e. are better able to exploit copying. Note in particular that shingling does the best at identifying tweets that are completely specified by a copy of another tweet (when the number of residual edges is zero).

Let the residual edges per link be the number of residual edges in the encoding divided by the number of edges in the graph. If there is no copying in the encoding, then this number is 1; if the encoding is achieved by copying alone, then it is 0. Figure 3(b) shows that for our dataset, there is a linear relationship between the bits per link (y-axis) and the residual edges per link (x-axis), across all orderings.



(b) Bits per link vs. residual edges per link.

Fig. 3. (a) Number of tweets (in thousands, y-axis) with a certain number of residual edges (x-axis), for shingle (blue circles), hashtag (magenta x's), natural (cyan squares), and random (red crosses) orderings. (b) Bits per link versus the number of residual edges per link, for different orderings. The dotted line corresponds to the equation, $y = 20.672x + 2.628$.

We speculate that the different character of our plots in Figure (2) compared to those observed by Chierichetti et al has to do with differences in the data that we chose to compress, namely, that our Twitter graph represents natural language data (recall that it is a bipartite graph of tweets and words in tweets) rather than a social network or Web graph. In particular, recall that the number of bits used to encode each residual edge for a particular node is set by the size of the largest gap. It seems that this number is not particularly affected by the ordering, which does not seem unreasonable for sets of residual edges that correspond to sets of co-occurring words.

7. FUTURE DIRECTIONS Our study has been fairly limited in scope. For our semantic ordering heuristic, we only considered semantic

features of nodes; it would be interesting to also consider semantics on edges. Also, we only considered the most simple methods for using and combining semantic features for producing orderings based on lexicographic (or numeric) sorts. It may be worth exploring more sophisticated methods for combining multiple semantic features to produce an ordering or clustering of nodes.

These real world graphs typically live in databases, where compression must (1) accompany efficient partitioning across files on disk, and possibly sharding across multiple servers for very large graphs, and (2) be balanced with efficient access methods for query execution, as there are usually trade-offs between the two. It is worth noting that the lexicographic and/or numeric sorts employed by our simple semantic ordering heuristics may correspond to a natural layout of data within a database and/or useful database indexes (which can be efficiently compressed because they are sorted). As mentioned earlier,

it is important

3“Retweeting” – when a user copies another user’s tweet – is apparently popular among Twitter users.
7

Count (thousands)

Bits per link

to evaluate graph compression schemes not only with respect to storage performance, but also in the context of real systems with data access needs.

We have only studied one particular dataset, a corpus of short messages that we treated as “bags of words” and represented as a bipartite graph. While our dataset is not standard in the literature for evaluating graph compression performance, it allowed us to demonstrate encouraging results for our semantic ordering heuristic. As additional side effects, this dataset also gave us an excuse to consider compression for bipartite graphs and in particular, graphs associated with natural language.

8. CONCLUSIONS The Web graph’s property of lexicographic locality allows an efficient ordering heuristic – lexicographically

sorting the nodes by URL, i.e. node “name” – that can be exploited by a scheme like BV compression. We essentially extended the notion of lexicographic locality to include any semantic features on the graph, such as additional structured attributes associated with nodes. This idea is especially useful for graphs whose node names alone do not provide a natural ordering, if the nodes are associated with other semantic information that corresponds to some notion of locality, such as social networks, whose nodes correspond to individuals who each have a structured user profile.

9. ACKNOWLEDGEMENTS The author would like to thank M. Seltzer for suggesting edge labels as a heuristic for clustering nodes in a graph, V. Kanade for providing references about community detection and influencing the language in § 4.2, and M. Mitzenmacher for many fun and relevant readings, lectures, as signments, comments, etc.

REFERENCES Bol di, P., Rosa, M., Santini, M., and Vigna, S. 2010. Layered label propagation: A multiresolution ordered ordering

for compressing social networks. CoRR abs/1011.5425. Bol di, P. and Vigna, S. 2004. The webgraph framework: compression techniques. In WWW ’04: Proceedings of the 13th

international conference on World Wide Web. ACM, New York, NY, USA, 595–602. Broder, A. Z., Cha rkar, M., Frieze, A. M., and Mitzenmacher, M. 1998. Min-wise Independent Permutations. Journal

of Computer and System Sciences 60, 3 27–336. Chierichetti, F., Kumar, R., Lattanzi, S., Mitzenmacher, M., Panconesi, A., and Raghavan, P. 2009. On compressing

social networks. In KDD ’09: Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, New York, NY, USA, 219–228.

Feder, T. and Motwani, R. 1995. Clique partitions, graph compression and speeding-up algorithms. J. Comput. Syst. Sci. 51, 26 1–272.

Flake, G. W., Lawrence, S., and Giles, C. L. 2000. Efficient identification of web communities. In Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining. KDD ’00. ACM, New York, NY, USA, 150–160.

Ino, H., Kudo, M., and Nakamura, A. 2005. Partitioning of web graphs by community topology. In Proceedings of the 14th international conference on World Wide Web. WWW ’05. ACM, New York, NY, USA, 661–669.

Neumann, T. and Weikum, G. 2008. Rdf-3x: a risc-style engine for rdf. Proc. VLDB Endow. 1, 1, 647–659. Newman, M. E. J. 2003.

Fast algorithm for detecting community structure in networks. Newman, M. E. J. 2006. Modularity and community structure in networks. Proceedings of the National Academy of Sciences

103, 23 (June), 8577–8582. Yang, J. and Leskovec, J. 2010. 47.6 million twitter tweets. <http://snap.stanford.edu/data/twitter7.html>.